

1 We thank the reviewers for their overall positive feedback.

2 The main concern raised by **Reviewer #2** is that Assumption 3, which requires that the positive examples are sampled
3 within a margin γ from the boundaries of the set K_n , makes the overall results too weak. Reviewer #2 rightly noted
4 that for the distributions presented in the paper, γ should decay exponentially with n , and this may seem to be a strong
5 requirement. While this is a valid concern, we stress that we did not use this assumption at all for the results on failure
6 of gradient-descent. Since having a margin γ can only help the optimization, dropping this assumption simply makes
7 the optimization harder, hence these results still hold. Similarly, other negative results in the paper, namely - the
8 inability of shallow networks to express fractal distributions, hold without Assumption 3: this assumption only makes
9 the approximation problem easier.

10 In fact, the only results that rely on Assumption 3 is Theorem 1 and its corollaries, which give positive results, stating
11 that fractal distributions can be efficiently expressed by deep networks. While the existence of a margin simplifies
12 the construction made in the proof of this theorem, we can prove this theorem even **without** Assumption 3. We give
13 a sketch of such proof below. **To summarize, in order to answer the concern of Reviewer #2 we will completely**
14 **remove Assumption 3 from the final version, and adjust all the theorems accordingly.**

15 Following the suggestion of **Reviewer #3**, we ran an ex-
16 periment on the Vicsek distribution of depth 6, where
17 the examples are concentrated on the “fine” details of
18 the fractal. Such distribution is hard to approximate by
19 a shallow network, as shown in our theoretical analysis.
20 We trained networks of various depth and width on this
21 distribution (as in the experiments described in the origi-
22 nal submission). The results are shown in Figure 1. As
23 could be seen clearly, unlike distributions with “coarse”
24 approximation curve (shown in the original submission),
25 in this case the benefit of depth is not noticeable, and all
26 architectures achieve an accuracy of slightly more than
27 0.5 (i.e., chance level performance).

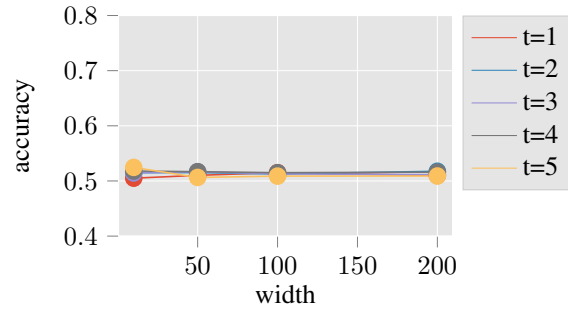


Figure 1: Performance on the “fine” Vicsek distribution.

28 We will additionally fix other minor issues raised by the
29 reviewers in the final version.

30 Proof Sketch of Theorem 1 without Assumption 3

31 **Lemma 2** (without Assumption 3, standard construction of a ReLU network) There exists a neural-network with two
32 hidden-layers such that $\mathcal{N}_{\mathbf{W},\mathbf{B}}(\mathbf{x}) < 0$ for $\mathbf{x} \notin [0, 1]^d$, and $\mathcal{N}_{\mathbf{W},\mathbf{B}}(\mathbf{x}) \geq 0$ for $\mathbf{x} \in [0, 1]^d$.

33 **Lemma 3** (without Assumption 3) There exists a neural-network of width $\max\{dr, 3d\}$ with two hidden-layers ($k =$
34 $3dr, t = 3$) such that for any n we have: $\mathcal{N}_{\mathbf{W},\mathbf{B}}(K_n) \subseteq K_{n-1}$ and $\mathcal{N}_{\mathbf{W},\mathbf{B}}(K_1 \setminus K_n) \subseteq \mathcal{X} \setminus K_{n-1}$

35 **Proof** Simple modification to the proof of Lemma 3 in the original submission. ■

37 **Lemma 4** (without Assumption 3) There exists a neural-network of width $2dr$ with two hidden-layers ($k = 2dr, t = 3$)
38 such that for any n we have: $\mathcal{N}_{\mathbf{W},\mathbf{B}}(\mathcal{X} \setminus K_1) < 0$ and $\mathcal{N}_{\mathbf{W},\mathbf{B}}(K_1) \geq 0$.

39 **Proof** Using Lemma 2, and following the same proof of Lemma 4 in the original submission. ■

40
41 **Proof** of Theorem 1 (without Assumption 3). We follow a proof similar to the one given in the original submission.
42 Instead of the original definition of h , we define $h(\mathbf{x}) = [g(\mathbf{x}_{1\dots d}), \mathbf{x}_{d+1} - \sigma(\mathbf{x}_{d+1} - \tilde{g}(\mathbf{x}_{1\dots d}))]$. Then, constructing
43 H as in the original proof satisfies that $H(\mathbf{x})_{d+1} < 0$ if and only if $\mathbf{x} \notin K_n$: if the $d + 1$ coordinate of some layer
44 becomes negative, it stays negative throughout the network (since the $d + 1$ coordinate of each layer is just the minimum
45 of the $d + 1$ coordinates of previous layers). Therefore, a network that outputs $H(\mathbf{x})_{d+1}$ achieves the required. ■

46